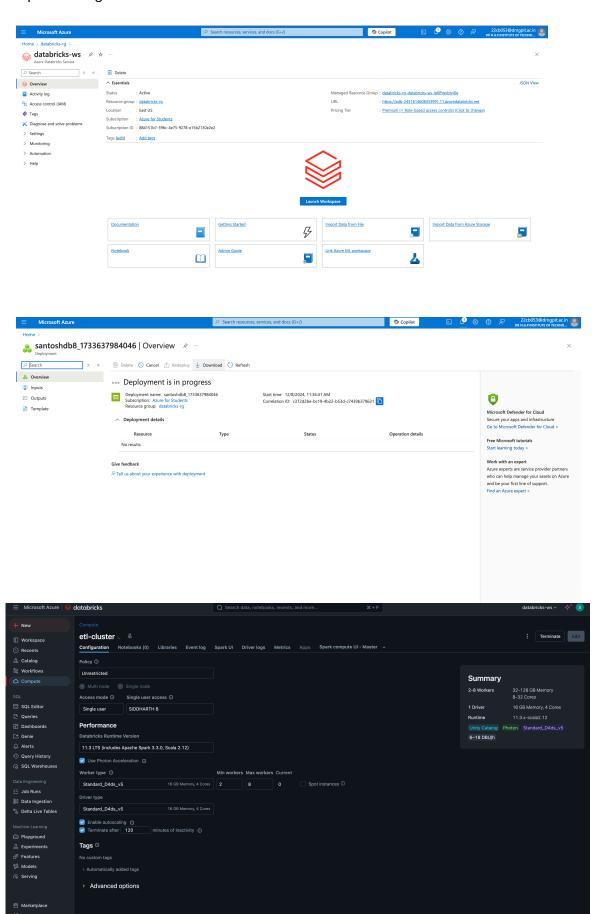# Step1 : Setting database containers and Azure data bricks

Step 2: Run the python code in Azure blob notebook (Load large dataset, preprocess, Implement ml ):